

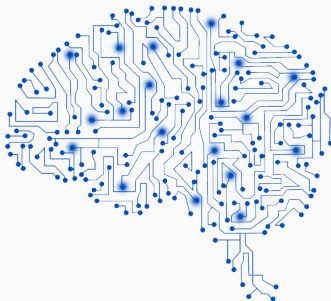
Apprentissage pour l'image

Machine learning for image processing

Presentation of TP1

Emile Pierret

Jeudi 9 mars 2023



At the end of the course :

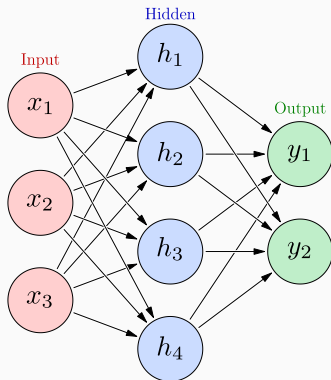
- Know what is a CNN (Convolutional Neural Network)
- Implement the training of a CNN for classification with Pytorch

For this session :

- Train the last layer of a multilayer network for classification with numpy.
- Apply a Stochastic gradient.
- For the moment, no images, no convolutions.

Reminder - What is a multilayer network ?

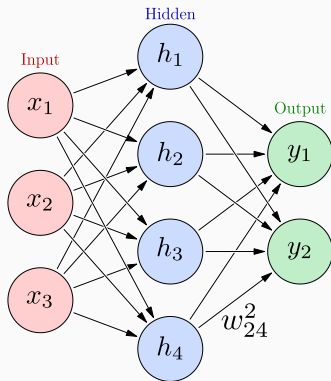
,containsverbatim]



- Inter-connection of several artificial neurons (also called nodes or units).
- Each level in the graph is called a layer:
 - Input layer,
 - Hidden layer(s),
 - Output layer.
- Each neuron in the hidden layers acts as a classifier / feature detector.
- Feedforward NN (no cycle)
 - first and simplest type of NN,
 - information moves in one direction.
- Recurrent NN (with cycle)
 - used for time sequences,
 - such as speech-recognition.

Reminder - What is a multilayer network ?

,containsverbatim]



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

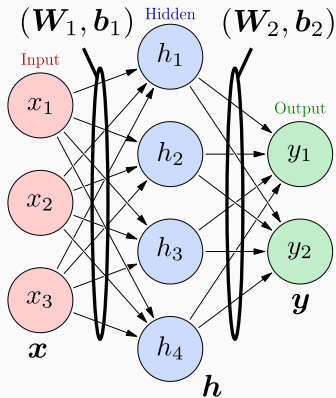
$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

Reminder - What is a multilayer network ?

,containsverbatim]



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

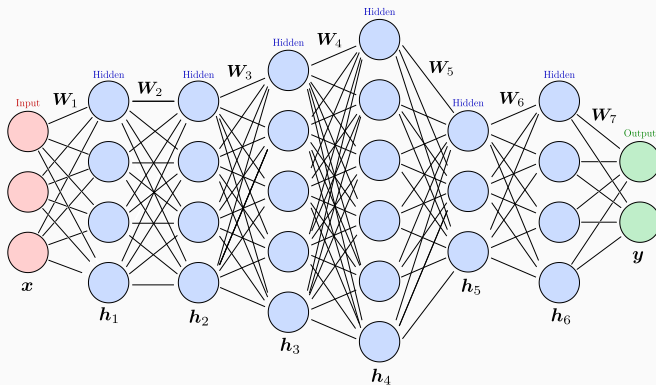
$$\mathbf{h} = g_1 (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

$$\mathbf{y} = g_2 (\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

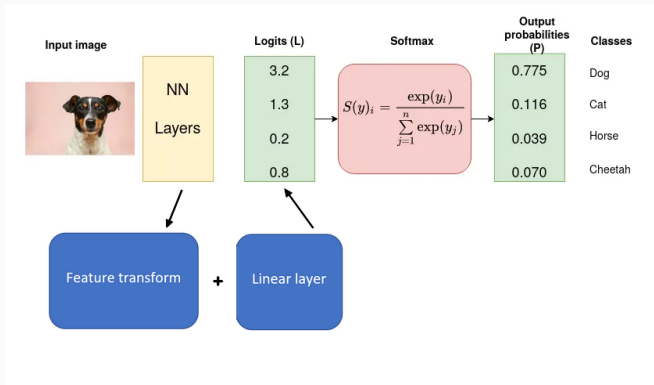
Artificial neural network / Multilayer perceptron



How to train a network ?

- ❶ Consider a training data set
- ❷ Consider a loss $L(W)$
- ❸ Minimize L applying a stochastic gradient

In real-life



In real-life, all the network have to be trained. The feature transform is trained with its linear separation.

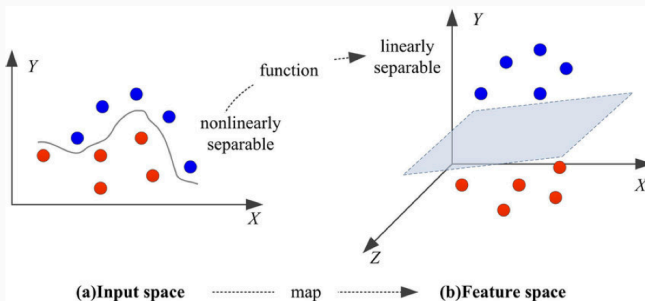
For this session, we will just train the last layer. We suppose that we have a good feature transform to have linearly separable data.

Feature transform

- We apply a feature transform $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^D$ to each \mathbf{x}_n :

$$\varphi_n = \varphi(\mathbf{x}_n), \quad n = 1, \dots, N.$$

- Depending on context it allows to increase ($D > p$) or decrease ($D < p$) the dimension in a way to favor class discrimination (e.g. PCA...).
- This is a non linear map that should make the classes linearly separable.



Recall of the notations

- K classes C_1, \dots, C_K
- (\mathbf{x}_n, t_n) the dataset such that $\mathbf{x}_n \in C_{t_n}$, $\mathbf{x}_n \in \mathbb{R}^d$
- $(\boldsymbol{\varphi}_n, t_n)$ the dataset after feature transform. $\boldsymbol{\varphi}_n \in \mathbb{R}^D$
- We suppose that $P(C_k|\boldsymbol{\varphi}) = \frac{\exp(\mathbf{w}_k^T \boldsymbol{\varphi})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \boldsymbol{\varphi})}$
- $\mathbf{y}_k(\boldsymbol{\varphi}_n) = \frac{\exp(\mathbf{w}_k^T \boldsymbol{\varphi}_n)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \boldsymbol{\varphi}_n)}$

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

or minimize

$$-\log \left(\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right) \right) = -\sum_{n=1}^N w_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

or minimize

$$-\log \left(\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right) \right) = -\sum_{n=1}^N \mathbf{w}_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- How to minimize it ?

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

or minimize

$$-\log \left(\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right) \right) = -\sum_{n=1}^N w_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- How to minimize it ?
- With the gradient descent !

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

or minimize

$$-\log \left(\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right) \right) = -\sum_{n=1}^N \mathbf{w}_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- How to minimize it ?
- With the gradient descent !
- What do we need ?

The feature transform being done we suppose that :

$$P(C_k|\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}$$

How to learn the weights $(\mathbf{w}_j)_{1 \leq j \leq K}$?

- The loss is the cross entropy, deriving from ?
- The maximization of the log-likelihood !
- More precisely, we would like to maximize :

$$\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right)$$

or minimize

$$-\log \left(\mathbb{P} \left(\bigcap_{n=1}^N \varphi_n \in C_{t_n} \mid W \right) \right) = -\sum_{n=1}^N \mathbf{w}_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- How to minimize it ?
- With the gradient descent !
- What do we need ?
- The gradient !

Computation of the gradient

$$L(\mathbf{W}) = - \sum_{n=1}^N w_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- Linear part: Partial gradient with respect to column w_ℓ , $\ell \in \{1, \dots, K\}$:

$$\begin{aligned} \nabla_{w_{t_n}} \left[w_{t_n}^T \varphi_n \right] &= \varphi_n \\ \nabla_{w_\ell} \left[w_{t_n}^T \varphi_n \right] &= 0 \text{ if } \ell \neq t_n \end{aligned}$$

- Partial gradient of $\nabla_{w_\ell} \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$?

$$\begin{aligned} \nabla_{w_\ell} \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right) &= \nabla_{w_\ell} \log \left(\exp(\mathbf{w}_\ell^T \varphi_n) + \text{constant} \right) \\ &=? \end{aligned}$$

A point on the gradients (reminder of the course "optimization"?)

Recall that for $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$, for $h \in \mathbb{R}^n$,

$$d(g \circ f)(x)(h) = dg(f(x))(df(x)(h)).$$

Here,

$$d(g \circ f)(x)(h_1) = \nabla(g \circ f)(x)^T h_1 \quad \text{for } h_1 \in \mathbb{R}^n$$

$$dg(x)(h_2) = g'(x)h_2 \quad \text{for } h_2 \in \mathbb{R}$$

$$df(x)(h_3) = \nabla f(x)^T h_3 \quad \text{for } h_3 \in \mathbb{R}^n$$

Consequently,

$$\nabla(g \circ f)(x)^T h = g'(f(x))\nabla f(x)^T h$$

So,

$$\nabla(g \circ f)(x) = g'(f(x))\nabla f(x)$$

Multivariate logistic regression

Gradient of log-likelihood:

Recall that for $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$

$$\nabla(g \circ f)(x) = g'(f(x))\nabla f(x).$$

Here,

$$g(t) = \log(\exp(t) + c) \quad g'(t) = \frac{\exp(t)}{\exp(t) + c}$$

$$f(\mathbf{w}_\ell) = \mathbf{w}_\ell^T \varphi_n \quad \nabla f(\mathbf{w}_\ell) = \varphi_n.$$

So,

$$\nabla_{\mathbf{w}_\ell} \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right) = \frac{\exp(\mathbf{w}_\ell^T \varphi_n)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n)} \varphi_n = \mathbf{y}_\ell(\varphi_n) \varphi_n$$

since

$$\mathbf{y}_k(\varphi) = \frac{\exp(\mathbf{w}_k^T \varphi)}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi)}, \quad k = 1, \dots, K.$$

Gradient of log-likelihood:

$$L(\mathbf{W}) = - \sum_{n=1}^N w_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(\mathbf{w}_j^T \varphi_n) \right)$$

- For each column $w_\ell \in \mathbb{R}^D$ of \mathbf{W} , $\ell \in \{1, \dots, K\}$,

$$\nabla_{w_\ell} L(\mathbf{W}) = \sum_{n=1}^N (\delta_{\ell, t_n} - \mathbf{y}_{t_n}(\varphi_n)) \varphi_n \in \mathbb{R}^D$$

- OK with intuition ?

One word on the stochastic gradient descent

We would like to minimize w.r.t \mathbf{W} ,

$$\begin{aligned} L\left(\mathbf{W}, (\varphi_n, t_n)_{1 \leq n \leq N}\right) &= - \sum_{n=1}^N w_{t_n}^T \varphi_n + \log \left(\sum_{j=1}^K \exp(w_j^T \varphi_n) \right) \\ &= \sum_{n=1}^N \ell(\mathbf{W}, \varphi_n, t_n) \end{aligned}$$

Algorithm: (stochastic) gradient descent for $L(\mathbf{w})$

- Initialize \mathbf{W} randomly
- Repeat until convergence
 - For all (φ_n, t_n) , $1 \leq n \leq N$
 - For each \mathbf{w}_k , $1 \leq k \leq K$
Update: $\mathbf{w}_k \leftarrow \mathbf{w}_k - \gamma \nabla_{\mathbf{w}_k} \ell(\mathbf{W}, \varphi_n, t_n)$
- γ is called the learning rate.

Comparison with constant step gradient descent

Algorithm: Constant step gradient descent for $f(x)$

- Initialize x randomly
- Repeat until convergence
 - Update $x \leftarrow x - \gamma \nabla_x f(x)$
- Why "stochastic" previously ?

Comparison with constant step gradient descent

Algorithm: Constant step gradient descent for $f(x)$

- Initialize x randomly
- Repeat until convergence
 - Update $x \leftarrow x - \gamma \nabla_x f(x)$
- Why "stochastic" previously ?
- Because we suppose that the data follow a distribution \mathbb{P}_{data} and we would like to minimize

$$\mathbb{P}_{\text{data}}(\varphi \in C_k \mid W) = \mathbb{E}_{\mathbb{P}_{\text{data}}} [\mathbf{1}_{\varphi \in C_k} \mid W]$$

w.r.t W

- In our context, we suppose that our dataset are samples that represents the distribution of the distribution of the features.

Theoretical stochastic gradient descent

We would like to minimize :

$$L(\mathbf{W}) = \mathbb{E}_{\mathbf{u}}(\ell(\mathbf{W}, \mathbf{u}))$$

Algorithm: (stochastic) gradient descent for $L(\mathbf{W})$

- Initialize \mathbf{W} randomly
- Repeat until convergence
 - Sample $\mathbf{u} \sim \mathbb{P}_{\mathbf{u}}$ independent from the previous samples
 - Update : $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} \ell(\mathbf{W}, \mathbf{u})$
- In practice, γ is not constant.
- There exists theoretical results of convergence.

What is it necessary to code ?

- ① A toy dataset
- ② The gradient of the loss
- ③ The gradient descent
- ④ An evaluation of the outputs.